



HAL
open science

Optimal software-implemented Itoh–Tsujii inversion for F_{2^m}

Jeremy Maitin-Shepard

► **To cite this version:**

Jeremy Maitin-Shepard. Optimal software-implemented Itoh–Tsujii inversion for F_{2^m} . WCC2015 - 9th International Workshop on Coding and Cryptography 2015, Anne Canteaut, Gaëtan Leurent, Maria Naya-Plasencia, Apr 2015, Paris, France. hal-01275712

HAL Id: hal-01275712

<https://inria.hal.science/hal-01275712>

Submitted on 18 Feb 2016

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Optimal software-implemented Itoh–Tsujii inversion for \mathbb{F}_{2^m}

Jeremy Maitin-Shepard

University of California at Berkeley
jbms@cs.berkeley.edu

Abstract. Field inversion in \mathbb{F}_{2^m} dominates the cost of modern software implementations of certain elliptic curve cryptographic operations, such as point encoding/hashing into elliptic curves. [7, 6, 2] Itoh–Tsujii inversion using a polynomial basis and precomputed table-based multi-squaring has been demonstrated to be highly effective for software implementations [19, 14, 2], but the performance and memory use depend critically on the choice of addition chain and multi-squaring tables, which in prior work have been determined only by suboptimal ad-hoc methods and manual selection. We thoroughly investigated the performance/memory tradeoff for table-based linear transforms used for efficient multi-squaring. Based upon the results of that investigation, we devised a comprehensive cost model for Itoh–Tsujii inversion and a corresponding optimization procedure that is empirically fast and provably finds globally-optimal solutions. We tested this method on 8 binary fields commonly used for elliptic curve cryptography; our method found lower-cost solutions than the ad-hoc methods used previously, and for the first time enables a principled exploration of the time/memory tradeoff of inversion implementations.

1 Introduction

Field inversion in \mathbb{F}_{2^m} is a key primitive required by elliptic curve cryptographic operations over binary elliptic curves. On older CPU architectures, the total cost of field multiplications typically dominates the cost of all other operations due to the inefficiency of performing polynomial multiplication without native instruction support, with the cost of field inversion only about 3–8 times the cost of field multiplication. [8, 19] The fast carry-less multiplication in modern CPU architectures, such as the PCLMULDQ instruction supported by Intel’s Westmere and later processors, the ARMv7 VMULL.P8 instruction, and the ARMv8 PMUL/PMULL/PMULL2 instructions, reduces the cost of multiplication by a factor of 20 to 40, but does not significantly reduce the cost of inversion. Consequently, for applications such as point encoding/hashing into elliptic curves [7, 6, 2], field inversion can account for more than half of the total computational cost in a software implementation on modern hardware, particularly when using binary fields of large degree. A reduction in the cost of field inversion can, therefore, have a significant effect on the overall cost. For this reason, we investigated the optimal software implementation of field inversion for arbitrary binary fields \mathbb{F}_{2^m} . **Our contribution:** Based on the Itoh–Tsujii polynomial-basis method [19, 9], apparently the fastest method for binary field inversion in software, we developed an empirically-fast algorithm for searching over the space of field inversion implementations in order to provably minimize computational cost subject to an arbitrary bound on, or cost function for, memory use. This algorithm, based on a CPU architecture-specific cost model estimated from performance benchmarks, allows us to compute a time/memory Pareto frontier of inversion implementations for a given binary field and CPU architecture, and has allowed us to find lower-cost implementations for commonly-used binary fields than were obtained by the ad-hoc methods [19, 14] employed previously.

2 Itoh–Tsujii inversion

The inverse x^{-1} of an element $x \in \mathbb{F}_{2^m}^\times$ can be computed using the Itoh–Tsujii method [11] as follows: by Lagrange’s theorem, $x^{2^m-1} = 1$, which implies that $x^{2^m-2} = x^{-1}$. To perform this field exponentiation efficiently, an *addition chain* is used.

Definition 1. A sequence $\langle a_1, b_1 \rangle, \dots, \langle a_k, b_k \rangle$ is said to be an addition chain for p if $a_k = p$ and $\{b_i, a_i - b_i\} \subseteq \{a_j \mid j \leq i\} \cup \{1\}$.

Given a valid addition chain $\langle a_i, b_i \rangle_{i \in [1, k]}$ for $m-1$, the Itoh–Tsujii method iteratively computes

$$x^{2^{a_i}-1} = (x^{2^{a_i-b_i}-1})^{2^{b_i}} \cdot x^{2^{b_i}-1} \quad \text{for } i = 1, \dots, k; \quad (1)$$

$$x^{-1} = x^{2^m-2} = (x^{2^{m-1}-1})^2. \quad (2)$$

The two operations required by eq. (1) are the raising of the previously-computed term $x^{2^{a_i-b_i-1}}$ to the power 2^{b_i} , which is called multi-squaring [19] or m -squaring [1, 5], and the multiplication of the newly computed $(x^{2^{a_i-b_i-1}})^{2^{b_i}}$ term by the previously computed $x^{2^{b_i-1}}$ term. Equation (2) requires only a single squaring operation.

When using a polynomial basis [19, 9], multi-squaring to the power 2^b can be performed in one of two ways:

1. **directly** as a sequence of b squarings;
2. using a **lookup table**, as described in section 4, by taking advantage of the fact that squaring, and therefore multi-squaring to any fixed power, is linear in the binary coefficients. [5, 19, 14, 2]

This can be faster than b repeated squarings for sufficiently large values of b .

Note that any addition chain $\langle a_i, b_i \rangle_{i \in [1, k]}$ has a complementary addition chain $\langle a_i, a_i - b_i \rangle_{i \in [1, k]}$; while for some purposes no distinction need be made between these two chains, in the case of Itoh–Tsuji inversion the distinction is necessary because the two chains lead to a different sequence of operations.

3 Optimal addition chains

Prior work has demonstrated software implementations of Itoh–Tsuji inversion with multi-squaring tables to be highly efficient. [19, 14, 2] The cost depends critically, however, on the choice of addition chain and multi-squaring tables. The key problem we investigated is the optimization of the choice of addition chain and multi-squaring tables to minimize this cost, for arbitrary binary fields \mathbb{F}_{2^m} . In the case that only direct squaring is used (no precomputed multi-squaring tables), Itoh–Tsuji inversion using an addition chain of length k requires exactly k field multiplications and m squarings. Minimizing the cost is simply a matter of finding a minimum-length addition chain for $m - 1$, and does not depend on the actual cost of multiplication and squaring. (There are typically many minimum-length chains for a given value of m .) There is no known polynomial-time algorithm for finding minimum-length addition chains, but there are search algorithms based on various pruning criteria [20, 3] that are very effective in practice for reasonable values of m .¹

The use of precomputed multi-squaring tables, crucial to achieving good performance, as shown in section 7, greatly complicates the optimization problem: because we may use either direct squaring or table-based multi-squaring for each step of the addition chain, we must consider the combined cost of the multiplications, squarings, and table lookups required; minimizing the length of the chain only minimizes the multiplication cost, which is typically less than half of the total cost. We must also consider that the table lookup cost is not fixed but indirectly depends on the total number of multi-squaring tables used, as described in section 4. Furthermore, we must consider not only the computational cost but also the memory use of the precomputed tables; it may be desirable to trade off some performance for a reduced memory footprint. The optimal solution clearly *does* depend on the actual costs of multiplication, squaring, and table lookups, which depend on the CPU architecture and their particular implementation. Our measurement of the computational costs of these operations is described in section 5.

We devised a comprehensive architecture-specific cost model based on these measurements; we also developed an A* search procedure for efficiently optimizing this cost model, as described in section 6. As the prior work on pruning criteria for finding minimum-length addition chains is not directly applicable to our more complicated cost model, we defined a novel efficiently-computable search heuristic that does not sacrifice optimality under our cost model.

4 Table-based linear field operations

Several key operations for \mathbb{F}_{2^m} , such as squaring, multi-squaring, square root, and half-trace, are linear in the binary coefficients. For multi-squaring (useful for inversion) and half-trace, an implementation based on a lookup table can be significantly faster than direct computation. [5, 19, 14, 2] The coefficients are split into $\lceil m/B \rceil$ blocks of B bits, and a separate table of 2^B entries is pre-computed for each block position, using a total of $s_{m,B} = \lceil m/B \rceil \cdot 2^B \cdot \lceil m/W \rceil \cdot W/8$ bytes of memory, where W is the word size in bits. The linear transform can then be computed from the precomputed tables with $k = \lceil m/B \rceil \cdot \lceil m/W \rceil$ memory accesses and $k - 1$ XOR operations.

Reasonable values of B range from 4 to 16, with larger values reducing the number of memory accesses required per transform at the cost of greater memory use. Due to the nature of the

¹ As we are primarily concerned with elliptic curve cryptography, binary field sizes of interest range from about $m = 127$ to about $m = 571$.

cache/memory hierarchy, however, memory access latency and throughput depend significantly on the extent to which the target is cached; this effect can be sufficiently large to outweigh the reduction in memory accesses required with larger values of B . To optimize table use we must consider, therefore, the amount of cache memory expected to be available, which depends in particular on the total size of all lookup tables required concurrently.

5 Performance evaluation

In order to measure the costs of low-level field operations (namely multiplication, squaring, and multi-squaring), and also to evaluate field inversion implementations, we developed an optimized implementation of binary field operations as an open-source C++ library. [12] Through the use of C++ templates, we were able to write much of the code generically to support different field sizes without sacrificing runtime performance; only for modular reduction was a custom implementation required for each supported field. We incorporated existing fast x86/x86-64 polynomial multiplication, squaring, and modular reduction routines for $\mathbb{F}_{2^{163}}$, $\mathbb{F}_{2^{193}}$, $\mathbb{F}_{2^{233}}$, $\mathbb{F}_{2^{239}}$, $\mathbb{F}_{2^{283}}$, $\mathbb{F}_{2^{409}}$, $\mathbb{F}_{2^{571}}$ [4] and for $\mathbb{F}_{2^{127}}$ [14]. These fields are sufficient for implementing elliptic curve operations for all NIST-recommended binary elliptic curves [18] as well as several other SEC 2-recommended binary elliptic curves [15]. $\mathbb{F}_{2^{127}}$ is needed to implement the record-breaking GLS254 curve [14] (based on the quadratic extension field $\mathbb{F}_{2^{254}}$).

As our test platforms we used an Intel Westmere i7-970 3.2 GHz CPU (with 12 MiB L3 cache) and an Intel Haswell i7-4790K 4.0 GHz CPU (with 8 MiB L3 cache). Both of these processors support the PCLMUL instruction for carry-less multiplication, Westmere being the first Intel architecture to support it; on the much more recent Haswell architecture, where this instruction has significantly lower cost, alternative modular reduction routines based on it are used for $\mathbb{F}_{2^{163}}$, $\mathbb{F}_{2^{283}}$, and $\mathbb{F}_{2^{571}}$ for a modest gain in performance.[4] All code was compiled separately for each architecture using version 3.5 of the Clang compiler at the highest optimization level.

5.1 Table-based operation performance

Since table-based multi-squaring, and field inversion routines based on it, induce data-dependent memory accesses, we measured the aggregate execution time for a set of inputs guaranteed to induce a uniform memory access pattern (and then divided by the number of inputs), in order to obtain worst-case warm-cache estimates. Failure to do so would result in a large underestimate of execution time due to caching.

We also observed the performance characteristics of table operations to be significantly affected by the size of the virtual memory pages backing the tables; in particular, on the x86-64 test machines, both the base level performance and the scaling of execution times with increasing table size were significantly better with 2 MiB (huge) pages than with 4 KiB pages, as shown in fig. 1, due to the cost of translation lookaside buffer (TLB) misses. The Linux transparent huge page support (introduced in Linux version 2.6.38) results in some, but not all, memory regions being backed automatically by huge pages, depending on a number of factors including region alignment and physical memory fragmentation; when not taken into account, this significantly reduced the reliability of our performance measurements. For consistent performance, we therefore ensured that all lookup tables were backed by huge pages.

We estimated the cost of table-based multi-squaring for each block size $B \in \{4, 6, 8, 10, 12, 16\}$ and for varying values $n \in [1, 6]$ of the number of concurrent lookup tables in use by timing n successive table-based linear transforms (using n separate tables), and then dividing by n . Our implementation used a word size of $W = 128$ bits.

6 A* search procedure

Recall that in order to invert an element $x \in \mathbb{F}_{2^m}$, Itoh–Tsuji inversion requires an addition chain $\langle a_1, b_1 \rangle, \dots, \langle a_k, b_k \rangle$ where $\{b_i, a_i - b_i\} \subseteq \{a_j \mid j \leq i\} \cup \{1\}$, and $a_k = m - 1$. We modeled the cost based on the following assumptions:

1. Each step of the addition chain incurs a fixed cost of c_{mul} for the required multiplication.
2. There is a cost of $b_i \cdot c_{\text{sq}}$ for each direct multi-squaring to the power 2^{b_i} .
3. If table-based multi-squaring is used in place of direct multi-squaring, there is a cost $\tau_{B,n}$ that depends on the table block size B and the number of distinct multi-squaring tables $n = |T|$ (where T is the set of distinct values b_i for which a multi-squaring table is used).

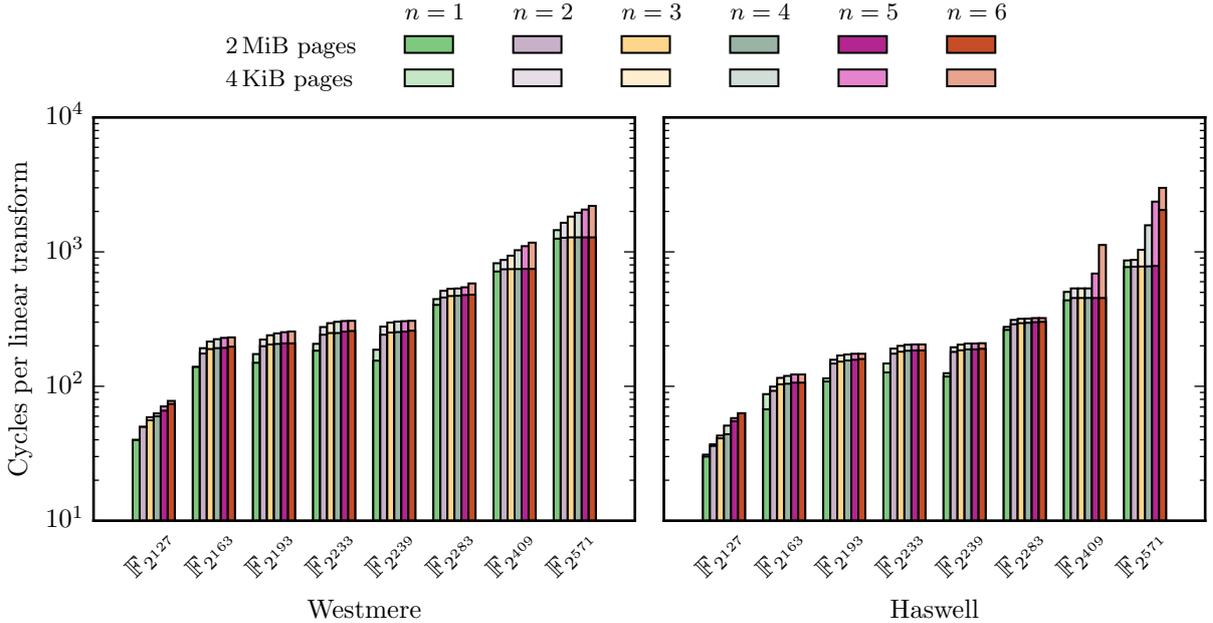


Fig. 1: Effect of virtual memory page size on table-based linear field operation (i.e. multi-square) performance, for varying values of n , the total number of lookup tables used concurrently. Results are shown for a fixed block size of $B = 8$ bits. The lighter-color bar extensions show the overhead introduced by the use of 4 KiB pages, as discussed in section 5.1.

We impose the minor restriction that the same block size B is used for all multi-squaring tables. In addition to the computational costs c_{mul} , c_{sq} , and $\tau_{B,n}$, we assume that there is a non-decreasing cost $c_{\text{mem}}(s)$ of consuming s bytes of memory with multi-squaring tables, which specifies a trade-off between computation time and memory use. Note that a hard memory constraint can be obtained by defining $c_{\text{mem}}(s)$ to be a step function to infinity.

6.1 Search formulation

To compute an optimal addition chain for $m-1$, we use as a search space the set S of valid addition chains $A = \langle a_i, b_i, o_i \rangle_{i \in [1,k]}$ summing to values $a_k \leq m-1$, where the extra variable $o_i \in \{\delta, \tau\}$ specifies whether direct (δ) or table-based (τ) multi-squaring is used. The successors of a state are the valid extensions of that addition chain. Without loss of generality, we consider only addition chains in which a_1, \dots, a_k are monotonically increasing.

Given an addition chain $A = \langle a_i, b_i, o_i \rangle_{i \in [1,k]}$, let $T_j(A) := \{b_i \mid i \in [1,j], o_i = \tau\}$ be the set of distinct multi-squaring tables used by the first j steps of the chain, and let $U_j^\tau(A) := \{i \mid i \in [1,j], o_i = \tau\}$ and $U_j^\delta(A) := \{i \mid i \in [1,j], o_i = \delta\}$ be the subsets of the first j steps for which table-based and direct multi-squaring are used, respectively. Restating our cost model formally, the cost of the j -step prefix $A_{1..j}$ of the chain A , depending on $n = |T_j(A)|$, is given by

$$\begin{aligned} \text{cost}(A_{1..j}) &:= \text{cost}_n(A_{1..j}) \\ &:= j \cdot c_{\text{mul}} + c_{\text{mem}}(n \cdot s_{m,B}) + |U_j^\tau(A)| \cdot \tau_{B,n} + \sum_{i \in U_j^\delta(A)} b_i \cdot c_{\text{sq}}. \end{aligned}$$

Recall that $s_{m,B}$ is the size in bytes of a single multi-squaring table with block size B , defined in section 4.

Table-based multi-squaring to the power 2^{a_i} is never advantageous relative to direct multi-squaring if $b_i \cdot c_{\text{sq}} \leq \tau_{B,|T|}$, and therefore we exclude such sub-optimal cases from the search space. Furthermore, in an optimal addition chain, if a multi-squaring table for the power 2^{b_i} is used at one step of the chain, i.e. $o_i = \tau$, then it must be the case that $c_{\text{sq}} \cdot b_i > \tau_{B,|T|}$, and thus there is no advantage in using direct multi-squaring for the same power at another step; therefore we also exclude this case from the search space.

6.2 Search algorithm

While the search space S is acyclic as defined, we can consider two equal-length addition chains $A = \langle a_i, b_i, o_i \rangle_{i \in [1, k]}$ and $A' = \langle a'_i, b'_i, o'_i \rangle_{i \in [1, k]}$ to be equivalent if the following conditions hold:

$$\begin{aligned} \{a_i \mid i \in [1, k]\} &= \{a'_i \mid i \in [1, k]\}, \\ \{(b_i, o_i) \mid i \in [1, k]\} &= \{(b'_i, o'_i) \mid i \in [1, k]\}, \\ |U_k^T(A)| &= |U_k^T(A')|. \end{aligned}$$

We use A^* graph search over the quotient set² of equivalence classes of S (each represented by the triplet $\langle \{a_i\}, \{(b_i, o_i)\}, |U_k^T| \rangle$) to find an optimal path;³ note that we still unambiguously specify *paths* through the quotient set as full addition chains $\langle a_1, b_1, o_1 \rangle, \dots, \langle a_k, b_k, o_k \rangle$. The initial state $\langle \emptyset, \emptyset, 0 \rangle$ corresponds to an empty addition chain.

Each search state is implicitly associated with a multi-squaring table block size B which is fixed after the initial state. The search priority queue is seeded with a separate initial state for each possible block size B .

For any given addition chain $A = \langle a_i, b_i, o_i \rangle_{i \in [1, k]} \in S$, A^* search requires a heuristic lower bound $h(A)$ on the cost to reach the goal. We specify this based on a *relaxation* of the problem:

Definition 2. A sequence $\langle a_i, b_i, o_i \rangle_{i \in [1, k]}$ is said to be a relaxed addition chain for p if

1. $a_k \geq p$, in place of the requirement that $a_k = p$;
2. $\max\{b_i, a_i - b_i\} \leq \max(\{a_j \mid j < i\} \cup \{1\})$ for all $i \in [1, k]$, in place of the requirement that $\{b_i, a_i - b_i\} \subseteq \{a_j \mid j < i\} \cup \{1\}$.

Definition 3. Let $\mathcal{R}_p(A)$ to be the set of relaxed addition chains $\langle a_1, b_1, o_1 \rangle, \dots, \langle a_k, b_k, o_k \rangle, \dots, \langle a_{k'}, b_{k'}, o_{k'} \rangle$ for p that extend $A = \langle a_i, b_i, o_i \rangle_{i \in [1, k]}$. We define the heuristic

$$h(A) := \min_{A' \in \mathcal{R}_{m-1}(A)} \text{cost}(A') - \text{cost}(A)$$

to be the minimum additional cost of a relaxed addition chain for $m - 1$ that extends A .

Since $h(A)$ is defined as the additional cost of the optimal solution to a relaxation of the problem, it is necessarily a consistent heuristic. This justifies our use of A^* graph search rather than tree search. [16, p. 95] We show that there is a very efficient algorithm for computing the $h(A)$, in contrast to the minimum *unrelaxed* addition chain cost, for which there is no polynomial-time algorithm.

Theorem 1. The heuristic cost $h(A)$ can be computed in $O(p \log p)$ time; furthermore, it depends only on a_k , $|T_k|$, $|U_k^T|$, and $\max T_k$, which respects the equivalence relation on addition chains.

Proof. See the extended version of the paper. [13]

7 Inversion performance evaluation

We implemented the A^* search procedure in Python and computed optimal Itoh–Tsuji inversion addition chains (for each binary field) for the Intel Westmere and Intel Haswell test platforms, estimating c_{sq} , c_{mul} , and $\tau_{B, n}$ separately for each platform as described in section 5. We varied c_{mem} in order to obtain all Pareto-optimal solutions (with respect to memory use and the model-predicted computational cost). We used a simple program to generate the C++ inversion code and appropriate multi-squaring tables for each addition chain computed by the search procedure, in order to evaluate the performance. The results are shown in fig. 2.

In addition to using the full cost model (referred to as the *variable table cost* (VT) model in the result figures and tables), we also computed optimal addition chains (using the same search algorithm) under several restricted models in order to evaluate the utility of each component of our model:

1. A *table-free* (TF) model in which $\tau_{B, n} = \infty$, such that only direct multi-squaring is used, and the search is reduced to merely finding a minimum-length addition chain. A table-free implementation, as used by Bluhm and Gueron [4], has the advantage of being constant time with respect to the input, which is necessary for the security of some applications though is often unnecessary in the hashing setting.

² It is straightforward to verify that this equivalence relation respects all search operations.

³ Concretely, we do not expand a state during the search if we have already expanded an equivalent state.

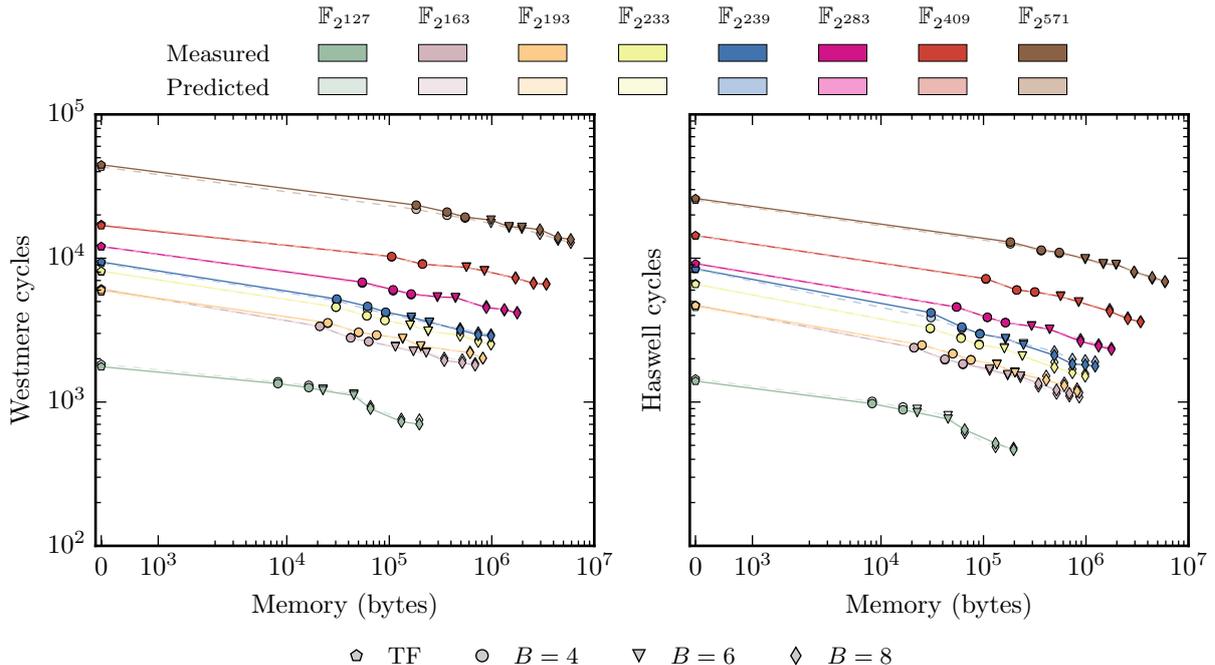


Fig. 2: Inversion time/memory Pareto frontiers, showing the trade-off between memory use and CPU cycles per field inversion. All Pareto-optimal points obtained by varying c_{mem} and jointly optimizing over both the addition chain and the multi-square table block size B are shown; the corresponding block size for each point is indicated by the marker. The execution times predicted by the full variable table cost model are shown with dashed lines in lighter shades. The points on the left with zero memory use correspond to table-free (TF) addition chains. No Pareto-optimal points using a block size $B \in \{10, 12, 16\}$ were found.

2. A *minimum-length* (ML) model in which table-based multi-squaring is used but not accounted for by the addition chain search. As in the table-free case, the search is reduced to merely finding a minimum-length addition chain, but multi-squaring tables are used for any value b_i in the chain for which $c_{\text{sq}} \cdot b_i > \tau$, for some threshold τ . Choosing $\tau = \tau_{B,1}$ roughly corresponds to the approach taken by prior work [14, 19] using multi-squaring tables.
3. A *single table cost* (ST) model in which we fix $\tau_{B,n} = \tau_{B,1}$, the cost of table-based multi-squaring estimated from a single-table benchmark.

Detailed results are given in appendix A.

8 Discussion

Reliable fine-grained performance measurement can be challenging due to compiler instruction reordering and optimization, processor branch prediction and out-of-order execution, and operating system jitter, among other factors, and these effects are magnified in our estimation of the cost of field inversion implementations as the sum of the costs of many individual operations. By carefully controlling for other sources of variance, including virtual page table size, the number of concurrent tables, and operation-specific benchmarking overhead, our cost model achieved a very low median error rate of 1.50% (fig. 3) in spite of these challenges.

Our A^* search procedure proved to be highly effective in computing optimal addition chains for any desired performance/memory tradeoff. Even with an unoptimized Python search implementation, our relaxation-based heuristic made the computation of complete Pareto frontiers (requiring many independent searches) take only a few CPU minutes (appendix B), a negligible amount of precomputation time for most elliptic curve cryptographic applications. In comparison, without a heuristic the search is completely intractable even for small field degrees.

The accuracy of our cost model allowed our search procedure to find addition chains that significantly reduced computational and memory costs compared to addition chains found by the

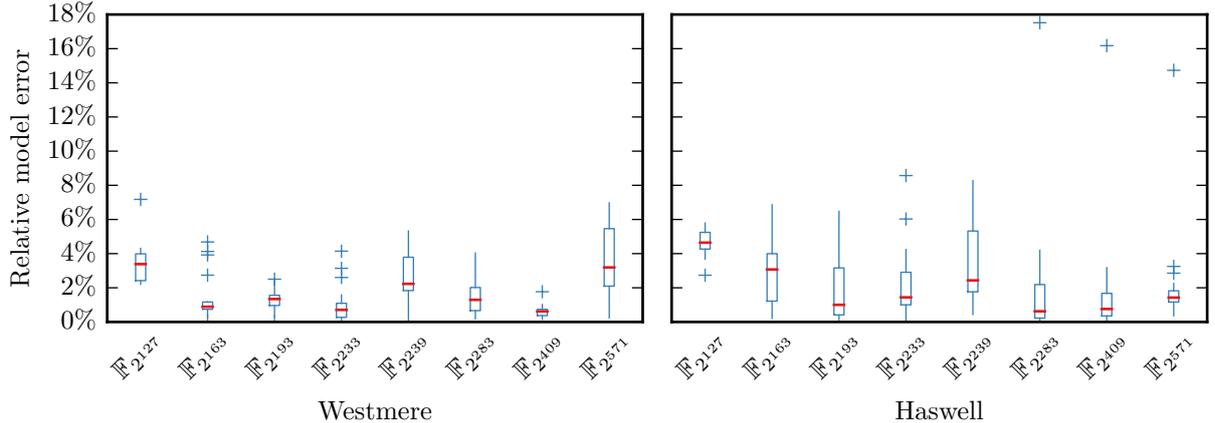


Fig. 3: Field inversion cost model relative errors (box plot). The error was calculated as the absolute value of the difference of the predicted cost and the measured cost, divided by the measured cost. The results were aggregated over all distinct optimal addition chains for all multi-square block sizes B and all values of c_{mem} . The overall median relative error was 1.50%.

ad-hoc criteria used by prior work (appendix A). Comparisons to simplified models confirm the utility of our full model: performance is improved by accounting for the cost of table-based multi-squaring in the addition chain optimization, and the performance is further improved by modeling the dependence of this cost on the number of concurrent multi-squaring tables in use. Compared to choosing the *best possible* set of multi-square tables for a fixed minimum-length addition chain, the joint model-based optimization over both addition chain and multi-square tables performed by our search procedure results in significantly better performance.

We imposed several limitations on our model, in the interest of keeping the complexity of the cost model and search procedure manageable: (a) inversion implementations are restricted to using a single block size B for all multi-squaring tables; and (b) other methods of multi-squaring are not considered, such as using a combination of one or more lookup tables and direct squaring. We expect, however, that any improvement in performance or memory use from lifting these restrictions would be small.

The use of multi-squaring tables introduces susceptibility to side channel attacks. However, blinding can be used to effectively prevent any information leak from memory access patterns, at the cost of just two additional field multiplications (and m random bits).

In the case that multiple independent field inversions must be computed, a speedup orthogonal to our investigation of optimal addition chains is also possible. Using Montgomery’s trick, n independent elements can be inverted simultaneously at the cost of just 1 field inversion and $3(n - 1)$ field multiplications. [17] Unless the batch size n is very large, however, it is still advantageous to reduce the cost of the single inversion required.

Alternative inversion methods exist, but are less efficient for software implementations on modern CPUs. Itoh–Tsujii inversion was originally proposed for use with a normal basis rather than polynomial basis representation of field elements. Under a normal basis $\{\beta^{2^i} \mid i \in [0, m - 1]\}$, multi-squaring is just a cyclic shift, with negligible computational cost. However, software-implemented multiplication is typically much slower with a normal basis than with a polynomial basis, [10, p. 72] and in practice this outweighs the speedup of multi-squaring. Inversion can also be performed using the Extended Euclidean Algorithm, but due to the large number of branching operations and shifts by arbitrary numbers of bits required, the performance is not competitive with polynomial basis Itoh–Tsujii inversion on modern CPU architectures. [19]

Bibliography

- [1] Ahmadi, O., Hankerson, D., Rodríguez-Henríquez, F.: Parallel formulations of scalar multiplication on Koblitz curves. *Journal of Universal Computer Science* 14(3), 481–504 (2008)
- [2] Aranha, D.F., Fouque, P.A., Qian, C., Tibouchi, M., Zapalowicz, J.C.: Binary elligator squared. *Cryptology ePrint Archive*, Report 2014/486 (2014), <http://eprint.iacr.org/>
- [3] Bahig, H.M.: Improved generation of minimal addition chains. *Computing* 78(2), 161–172 (2006)
- [4] Bluhm, M., Gueron, S.: Fast software implementation of binary elliptic curve cryptography. *IACR Cryptology ePrint Archive* 2013, 741 (2013)
- [5] Bos, J.W., Kleinjung, T., Niederhagen, R., Schwabe, P.: Ecc2k-130 on cell cpus. In: *Progress in Cryptology—AFRICACRYPT 2010*, pp. 225–242. Springer (2010)
- [6] Brown, D.R.L.: The encrypted elliptic curve hash. *IACR Cryptology ePrint Archive* 2008, 12 (2008)
- [7] Brown, D.R.L., Antipa, A., Campagna, M., Struik, R.: ECOH: the Elliptic Curve Only Hash. Submission to NIST (2008)
- [8] Fong, K., Hankerson, D., López, J., Menezes, A.: Field inversion and point halving revisited. *Computers, IEEE Transactions on* 53(8), 1047–1059 (2004)
- [9] Guajardo, J., Paar, C.: Itoh-Tsujii inversion in standard basis and its application in cryptography and codes. *Designs, Codes and Cryptography* 25(2), 207–216 (2002)
- [10] Hankerson, D., Vanstone, S., Menezes, A.J.: *Guide to elliptic curve cryptography*. Springer (2004)
- [11] Itoh, T., Tsujii, S.: A fast algorithm for computing multiplicative inverses in $\text{GF}(2^m)$ using normal bases. *Information and computation* 78(3), 171–177 (1988)
- [12] Maitin-Shepard, J.: C++ Elliptic Curve Multiset Hash library. <http://jeremys.com/ecmh> (2015)
- [13] Maitin-Shepard, J.: Optimal software-implemented Itoh-Tsujii inversion for $\text{GF}(2^m)$. *Cryptology ePrint Archive*, Report 2015/028 (2015), <http://eprint.iacr.org/>
- [14] Oliveira, T., López, J., Aranha, D.F., Rodríguez-Henríquez, F.: Two is the fastest prime: lambda coordinates for binary elliptic curves. *Journal of Cryptographic Engineering* 4(1), 3–17 (2014)
- [15] Research, C.: SEC 2: Recommended Elliptic Curve Domain Parameters. *Standards for Efficient Cryptography* (Sep 2000), http://www.secg.org/download/aid-386/sec2_final.pdf, version 1.0
- [16] Russell, S., Norvig, P.: *Artificial Intelligence: A Modern Approach*. Prentice Hall series in artificial intelligence, Prentice Hall (2010)
- [17] Shacham, H., Boneh, D.: Improving ssl handshake performance via batching. In: *Topics in Cryptology—CT-RSA 2001*, pp. 28–43. Springer (2001)
- [18] of Standards, N.I., Technology: FIPS 186-4: Digital Signature Standard (DSS), Federal Information Processing Standard (FIPS), publication 186-4. Tech. rep., Department of Commerce, Gaithersburg, MD, USA (Jul 2013), <http://nvlpubs.nist.gov/nistpubs/FIPS/NIST.FIPS.186-4.pdf>
- [19] Taverne, J., Faz-Hernández, A., Aranha, D.F., Rodríguez-Henríquez, F., Hankerson, D., López, J.: Software implementation of binary elliptic curves: impact of the carry-less multiplier on scalar multiplication. In: *Cryptographic Hardware and Embedded Systems—CHES 2011*, pp. 108–123. Springer (2011)
- [20] Thurber, E.G.: Efficient generation of minimal length addition chains. *SIAM Journal on Computing* 28(4), 1247–1263 (1999)

A Cost model comparison

For a fixed block size of $B = 8$ and a fixed cost $c_{\text{mem}} = 0$, the addition chains that resulted from each of the four cost model variants are compared in table 1. The optimal addition chains and multi-square table sets under each cost model for each CPU architecture are shown, along with the measured execution time, and the execution time predicted by our complete cost model.⁴

Relative to the table-free implementation, the use of multi-square tables with the same addition chain (the minimum-length model) provides a large improvement, reducing the computational cost by a factor of 2–4. Optimizing the addition chain as well, using the single table cost model, provided on average an additional moderate reduction in cost. The full variable table cost model further reduced the computational cost for some of the fields. While in some cases, such as for \mathbb{F}_{2409} , the same addition chain was optimal under multiple cost models, the full model never resulted in worse performance.

To compare the cost models under broader conditions, we generated inversion implementations using the same minimum length (ML) addition chains shown in table 1 for *each possible* threshold τ and block size B . We measured the performance of each implementation, and computed an *empirical* time/memory Pareto frontier. We compared this to the actual performance along the Pareto frontier obtained based on the *estimated* costs of optimal points under our complete model. The results are shown in fig. 4. Averaging over all Pareto-optimal minimum length implementations (and both the Westmere and Haswell CPU architectures), the optimal addition chains under our cost model required 15.3% less memory at the same performance level, or 6.7% less computation time at the same level of memory use. Using the minimum length criteria alone to select the addition chain essentially leaves the performance up to chance: different minimum-length addition chains, when used with multi-squaring tables, can have very different computational cost.

⁴ The predicted execution times shown are based on the complete cost model in all cases, even for entries corresponding to addition chosen computed based on one of the restricted models.

Table 1: \mathbb{F}_{2^m} inversion addition chain cost model comparison. The optimal addition chains computed for the Intel Westmere (W) and Intel Haswell (H) machines under the *table-free* (TF), *minimum length* (ML), *single table cost* (ST), and the full *variable table cost* (VT) models are shown, along with the associated actual execution time in CPU cycles and predicted execution time (predict.) under the full variable table cost model. A fixed block size of $B = 8$ and $c_{\text{mem}} = 0$ were used. The addition chains were optimized independently for each test machine and cost model, but in some cases the same optimal chain was found for multiple machines/cost models.

	Model	Addition chain (b_i)	Tables	Westmere Actual/Predict.	Haswell Actual/Predict.
$\mathbb{F}_{2^{127}}$	TF	1 1 1 3 7 7 21 21 63	–	1761 / 1848	1397 / 1460
	ML	1 1 1 3 7 7 21 21 63	7 21 63	721 / 774	515 / 490
	ST	1 1 2 5 1 11 21 21 63	5 11 21 63	726 / 783	485 / 495
	VT	1 2 2 4 1 11 21 21 63	11 21 63	701 / 764	465 / 489
$\mathbb{F}_{2^{163}}$	TF	1 1 2 5 10 20 1 40 81	–	6102 / 6160	4742 / 4732
	ML/ST/VT(H)	1 1 2 5 10 20 1 40 81	5 10 20 40 81	1853 / 1920	1206 / 1111
	VT(W)	1 2 2 4 10 20 1 41 81	10 20 41 81	1808 / 1890	
$\mathbb{F}_{2^{193}}$	TF	1 1 3 6 12 24 48 96	–	5996 / 5875	4673 / 4595
	ML/ST	1 1 3 6 12 24 48 96	6 12 24 48 96	2091 / 2105	1231 / 1291
	VT	1 2 4 4 12 24 48 96	12 24 48 96	2018 / 2044	1190 / 1258
$\mathbb{F}_{2^{233}}$	TF	1 1 1 3 4 7 11 29 58 116	–	8127 / 8071	6598 / 6696
	ML	1 1 1 3 4 7 11 29 58 116	7 11 29 58 116	2682 / 2744	1611 / 1681
	ST	1 1 1 3 7 1 15 29 58 116	7 15 29 58 116	2578 / 2654	1515 / 1601
	VT	1 1 3 2 6 1 15 29 58 116	15 29 58 116	2524 / 2584	1506 / 1600
$\mathbb{F}_{2^{239}}$	TF	1 1 1 3 7 14 28 7 56 119	–	9382 / 9271	8741 / 8402
	ML/ST/VT(H)	1 1 1 3 7 14 28 7 56 119	7 14 28 56 119	2912 / 3006	1783 / 1937
	VT(W)	1 1 3 2 6 14 28 56 14 126	14 28 56 126	2885 / 2977	
$\mathbb{F}_{2^{283}}$	TF	1 1 1 3 3 7 17 34 68 10 136	–	12071 / 12090	9174 / 9101
	ML(W)	1 1 1 3 3 7 17 34 68 10 136	17 34 68 136	4452 / 4540	
	ML(H)	1 1 1 3 3 7 17 34 68 10 136	10 17 34 68 136		2648 / 2650
	ST/VT	1 2 4 1 8 1 18 35 1 71 141	18 35 71 141	4145 / 4244	2317 / 2398
$\mathbb{F}_{2^{409}}$	TF	1 1 3 6 12 3 24 51 102 204	–	16993 / 17210	14375 / 14373
	ML/ST/VT	1 1 3 6 12 3 24 51 102 204	24 51 102 204	6676 / 6666	3631 / 3664
$\mathbb{F}_{2^{571}}$	TF	1 1 1 3 4 11 22 7 44 95 95 285	–	44828 / 43235	26063 / 25492
	ML	1 1 1 3 4 11 22 7 44 95 95 285	22 44 95 285	14568 / 13469	7301 / 7212
	ST/VT	1 1 2 2 7 7 14 35 70 5 145 285	35 70 145 285	13646 / 12922	6962 / 6883

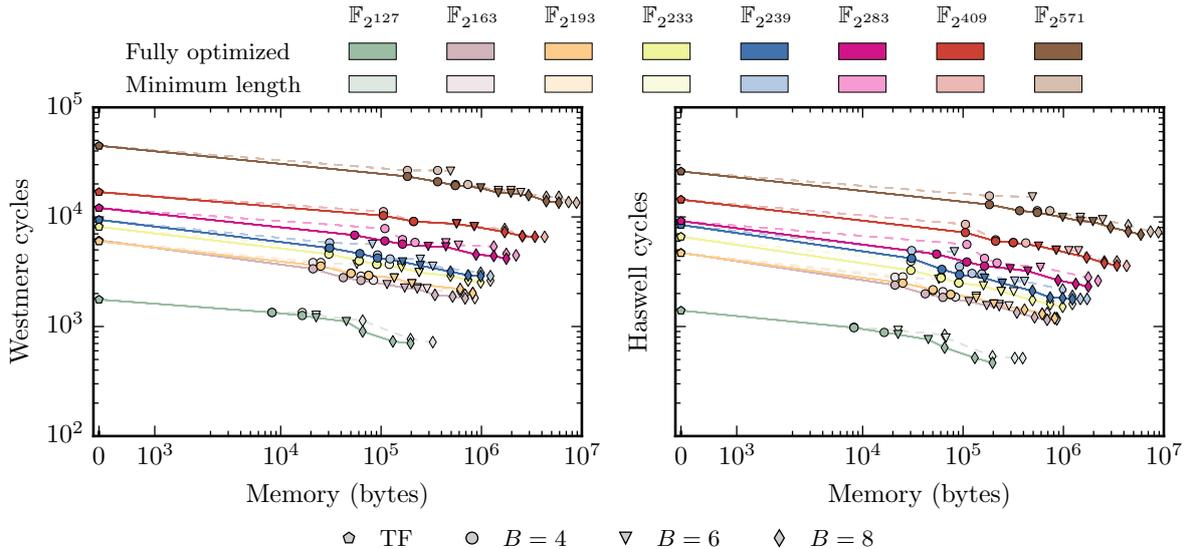


Fig. 4: Effect of cost model on inversion time/memory Pareto frontiers. For each field, a minimum-length addition chain was computed, and all inversion implementations obtained by varying the number of multi-square tables and the block size B were evaluated. The Pareto-optimal points (according to the actual timing results) for these minimum-length addition chain implementations are shown with dashed lines and in lighter colors, overlaid by the Pareto-optimal points obtained from our complete cost model by varying c_{mem} .

Table 2: A^* search performance for Pareto frontier computation. For each field and block size B , a set of Pareto-optimal addition chains were computed using several independent searches with successively lower memory bounds. The total number ($\#$) of searches performed, the total time required by all searches, the number of search states expanded (successors computed), and the number of heuristic calls are shown.

Field	Westmere model				Haswell model			
	#	Time (s)	States expanded	Heuristic calls Total / Uniq.	#	Time (s)	States expanded	Heuristic calls Total / Uniq.
\mathbb{F}_{2127}	11	1.04	4 881	48 144/ 3 518	15	1.21	5 805	56 196/ 4 200
\mathbb{F}_{2163}	20	0.56	2 640	25 285/ 4 872	21	0.51	2 317	22 456/ 5 015
\mathbb{F}_{2193}	19	0.10	598	3 803/ 1 311	19	0.12	691	4 774/ 1 520
\mathbb{F}_{2233}	19	2.09	8 877	98 210/ 9 067	18	8.48	25 895	408 592/ 10 143
\mathbb{F}_{2239}	19	51.35	161 499	2 401 801/ 55 121	19	83.85	227 070	3 871 342/ 48 086
\mathbb{F}_{2283}	18	13.10	51 094	588 551/ 24 589	17	9.39	36 464	425 197/ 13 807
\mathbb{F}_{2409}	15	0.58	2 883	24 642/ 2 932	15	1.24	5 301	55 089/ 5 204
\mathbb{F}_{2571}	16	115.97	352 675	4 861 047/ 57 307	16	157.53	471 921	6 478 492/ 51 543

B Search efficiency

Run times and other statistics for the search procedure itself, summed over all of the independent searches required to compute the time/memory Pareto frontiers for each field degree, are given in table 2. The run times measure the single-core CPU time of our Python implementation run on the Westmere test machine, which was used for all searches. As our focus was on algorithmic efficiency rather than low-level optimization of the search procedure, the reported run times should be seen only as a general guide. Even for a fixed field degree, the number of search states expanded, and therefore the run time, is highly dependent on the particular ratios of c_{sq} , $\tau_{B,n}$, and c_{mem} to c_{mul} . Given that a shortest path search over addition chains without the use of a heuristic would be completely intractable, it is clear from the run times and numbers of states expanded that our relaxed addition chain heuristic is highly effective. The large ratio of total heuristic calls to unique tuples $\langle a_k, |T_k|, |U_k^T|, \max T_k \rangle$ demonstrates the utility of memoizing heuristic values.